

# A Gentle Introduction to Dimension Reduction and Semiparametric Approach

Yanyuan Ma

University of South Carolina

December 11, 2015

# Outline

1. Introduction to Dimension Reduction Models
2. Inverse Regression Based Methods
3. Nonparametric Methods
4. Semiparametric Methods
5. Inference Issues
6. Determining the Structural Dimension
7. The Issue of High Dimension
8. Extensions

## Two Assumptions to Reduce Dimension

- ▶ One aspect of data complexity: number of covariates
- ▶ Assumption 1: Sparsity
  - ▶ Only a few of the covariates useful
  - ▶ Variable selection, often via penalization
- ▶ Assumption 2: Reducibility
  - ▶ Only a few linear combinations useful
  - ▶ (Sufficient) Dimension reduction

# Specific Dimension Reduction Models

- ▶ Notation:  $Y \in R$ ,  $\mathbf{x} \in R^{p \times 1}$ ,  $\beta \in R^{p \times d}$ ,  $p > d$ .
- ▶ Central space model
- ▶ Central mean space model
- ▶ Central moment space model
- ▶ Central variance space model
- ▶ Central median space model
- ▶ Central quantile space model

# Central Space Model

- ▶ Central space model assumption:  
Distribution of  $Y$  relates to covariates  $\mathbf{x}$  only through  $\beta^T \mathbf{x}$ .
- ▶ Mathematical form:

$$Y \perp\!\!\!\perp \mathbf{x} \mid \beta^T \mathbf{x}$$

or equivalently

$$\text{pr}(Y \leq y \mid \mathbf{x}) = \text{pr}(Y \leq y \mid \beta^T \mathbf{x}) \text{ for all } y.$$

- ▶ Central space  $\mathcal{S}_{Y|\mathbf{x}}$ : Span of the columns in  $\beta$ .
- ▶ Goal: Estimate  $\mathcal{S}_{Y|\mathbf{x}}$

# Central Mean Space Model

- ▶ Central mean space model assumption:  
Mean of  $Y$  relates to covariates  $\mathbf{x}$  only through  $\beta^T \mathbf{x}$ .
- ▶ Mathematical form:

$$E(Y | \mathbf{x}) = E(Y | \beta^T \mathbf{x})$$

or equivalently

$$Y = m(\beta^T \mathbf{x}) + \epsilon, \quad E(\epsilon | \mathbf{x}) = 0.$$

- ▶ Central mean space  $S_{E(Y|\mathbf{x})}$ : Span of the columns in  $\beta$ .
- ▶ Goal: Estimate  $S_{E(Y|\mathbf{x})}$

# Central Moment Space Model

- ▶ Central moment space model assumption:  
First  $k$  moments of  $Y$  relates to covariates  $\mathbf{x}$  only through  $\beta^T \mathbf{x}$ .
- ▶ Mathematical form:

$$E(Y^j | \mathbf{x}) = E(Y^j | \beta^T \mathbf{x}) \text{ for } j = 1, \dots, k.$$

or equivalently

$$Y^j = m_j(\beta^T \mathbf{x}) + \epsilon_j, \quad E(\epsilon_j | \mathbf{x}) = 0 \text{ for } j = 1, \dots, k.$$

- ▶ Central  $k$ -moment space  $S_{Y|\mathbf{x}}^{(k)}$ : Span of the columns in  $\beta$ .
- ▶ Goal: Estimate  $S_{Y|\mathbf{x}}^{(k)}$

# Central Variance Space Model

- ▶ Central variance space model assumption:  
Conditional variance of  $Y$  on  $\mathbf{x}$  is a function of  $\beta^T \mathbf{x}$  only.
- ▶ Mathematical form:

$$\text{var}(Y | \mathbf{x}) = E\{\text{var}(Y | \mathbf{x}) | \beta^T \mathbf{x}\}$$

- ▶ It is not equivalent to

$$\text{var}(Y | \mathbf{x}) = \text{var}(Y | \beta^T \mathbf{x})$$

- ▶ Central variance space  $S_{\text{var}(Y|\mathbf{x})}$ : Span of the columns in  $\beta$ .
- ▶ Goal: Estimate  $S_{\text{var}(Y|\mathbf{x})}$

# Estimation in Dimension Reduction Models

- ▶ Target: Finding the column space of  $\beta$  with the smallest  $d$ .
  - ▶ Estimating space, not parameters
  - ▶ Smallest space exists and unique
- ▶ Invariance: Let  $\mathbf{z} = \mathbf{A}^{-1}(\mathbf{x} - \mu)$ , then

$$\begin{aligned}\text{pr}(Y | \beta^T \mathbf{x}) &= \text{pr}\{Y | (\mathbf{A}\beta)^T \mathbf{z}\}, \\ \text{pr}(Y^j | \beta^T \mathbf{x}) &= \text{pr}\{Y^j | (\mathbf{A}\beta)^T \mathbf{z}\}, \\ \text{var}(Y | \beta^T \mathbf{x}) &= \text{var}\{Y | (\mathbf{A}\beta)^T \mathbf{z}\},\end{aligned}$$

which leads to

$$\begin{aligned}\mathcal{S}_{Y|\mathbf{x}} &= \mathbf{A}^{-1} \mathcal{S}_{Y|\mathbf{z}}, \quad \mathcal{S}_{E(Y|\mathbf{x})} = \mathbf{A}^{-1} \mathcal{S}_{E(Y|\mathbf{z})}, \\ \mathcal{S}_{Y|\mathbf{x}}^{(k)} &= \mathbf{A}^{-1} \mathcal{S}_{Y|\mathbf{z}}^{(k)}, \quad \mathcal{S}_{\text{var}(Y|\mathbf{x})} = \mathbf{A}^{-1} \mathcal{S}_{\text{var}(Y|\mathbf{z})}\end{aligned}$$

- ▶ Simplification:  $E(\mathbf{x}) = \mathbf{0}$ ,  $\text{cov}(\mathbf{x}) = \mathbf{I}_p$ .

# Three Classes of Estimation Approaches

1. Inverse regression based methods.
  - ▶  $\mathcal{S}_{Y|\mathbf{x}}$ : SIR, SAVE, DR
  - ▶  $\mathcal{S}_{E(Y|\mathbf{x})}$ : OLS, PHD
2. Nonparametric estimation:
  - ▶  $\mathcal{S}_{Y|\mathbf{x}}$ : dMAVE
  - ▶  $\mathcal{S}_{E(Y|\mathbf{x})}$ : MAVE
3. Semiparametric estimation:
  - ▶  $\mathcal{S}_{Y|\mathbf{x}}$ : Semi-SIR, Semi-DR, Semi-SAVE, Efficient
  - ▶  $\mathcal{S}_{E(Y|\mathbf{x})}$ : Semi-PHD, Local Efficient, Efficient

# Inverse Regression Based Methods

- ▶ Central idea: Invert the position between  $Y$  and  $\mathbf{x}$
- ▶ Requirements:
  - ▶ Linearity condition

$$E(\mathbf{x} | \boldsymbol{\beta}^T \mathbf{x}) = \mathbf{P}\mathbf{x} = \boldsymbol{\beta}(\boldsymbol{\beta}^T \boldsymbol{\beta})^{-1} \boldsymbol{\beta}^T \mathbf{x}$$

- ▶ Constant variance condition

$$\text{cov}(\mathbf{x} | \boldsymbol{\beta}^T \mathbf{x}) = \mathbf{Q} = \mathbf{I}_p - \mathbf{P}$$

- ▶ Linearity/constant variance condition need to hold at true  $\boldsymbol{\beta}$
- ▶ Linearity condition approximately holds for large  $p$ , fixed  $d$ .
- ▶ Hard to check. If strengthened to any  $\boldsymbol{\beta}$ : elliptical/normal distribution

# Sliced Inverse Regression (SIR)

- ▶ SIR estimate of  $S_{Y|\mathbf{x}}$ : non-zero eigen-space of  $\Lambda_{\text{sir}} = \text{cov}\{E(\mathbf{x} | Y)\}$ .
- ▶ Procedure:
  - ▶ Estimate  $\hat{E}(\mathbf{x} | Y)$ :  $p$  univariate problems, nonparametric regression or average in slices
  - ▶ Form sample variance-covariance of  $\hat{E}(\mathbf{x} | Y)$ .
  - ▶ Extremely simple procedure!
- ▶ SIR requirement: critically relies on linearity condition.

## Sliced Average Variance Estimation (SAVE)

- ▶ SAVE estimate of  $\mathcal{S}_{Y|\mathbf{x}}$ : non-zero eigen-space of
$$\Lambda_{\text{save}} = E \left[ (\mathbf{I}_p - \widehat{\text{cov}}(\mathbf{x} | Y))^2 \right].$$
- ▶ Procedure:
  - ▶ Estimate  $\widehat{\text{cov}}(\mathbf{x} | Y)$ :  $O(p^2)$  univariate problems, nonparametric regression or sample variance-covariance in slices
  - ▶ Form sample average of  $\{\mathbf{I}_p - \widehat{\text{cov}}(\mathbf{x} | Y)\}^2$ .
  - ▶ Computationally simple
- ▶ SAVE requirement: critically relies on linearity condition and constant variance condition.

# SIR or SAVE?

- ▶ SAVE more comprehensive
- ▶ SIR more efficient
- ▶ Convex combination of SIR and SAVE

$$\alpha \Lambda_{\text{sir}} + (1 - \alpha) \Lambda_{\text{save}}, \quad 0 \leq \alpha \leq 1$$

- ▶ More comprehensive than SIR, more efficient than SAVE.

# Directional Regression (DR)

- ▶ DR estimate of  $\mathcal{S}_{Y|x}$ : non-zero eigen-space of
$$\Lambda_{\text{dr}} = E[\{2\mathbf{I}_p - \mathbf{A}(Y, \tilde{Y})\}^2],$$
$$\mathbf{A}(Y, \tilde{Y}) = E\{(\mathbf{x} - \tilde{\mathbf{x}})(\mathbf{x} - \tilde{\mathbf{x}})^T \mid Y, \tilde{Y}\}$$
$$(\tilde{\mathbf{x}}, \tilde{Y}): \text{an independent copy of } (\mathbf{x}, Y).$$
- ▶ Procedure:
  - ▶ Estimate  $\hat{E}\{(\mathbf{x} - \tilde{\mathbf{x}})(\mathbf{x} - \tilde{\mathbf{x}})^T \mid Y, \tilde{Y}\}$ :  $O(p^2)$  univariate problems, nonparametric regression or sample average in slices
  - ▶ Form sample average of  $\{2\mathbf{I}_p - \hat{\mathbf{A}}(Y, \tilde{Y})\}^2$ .
  - ▶ Improve both SIR and SAVE
- ▶ DR requirement: critically relies on linearity condition and constant variance condition.

# Ordinary Least Squares (OLS)

- ▶ OLS estimate of  $\mathcal{S}_{E(Y|\mathbf{x})}$ :
  - ▶ Perform ordinary least square on model  $Y = \boldsymbol{\beta}^T \mathbf{X} + \epsilon$
  - ▶ The column space of  $\hat{\boldsymbol{\beta}}$  belongs to  $\mathcal{S}_{E(Y|\mathbf{x})}$ .
- ▶ Finds one direction, not all directions of  $\mathcal{S}_{E(Y|\mathbf{x})}$ .
- ▶ OLS requirement: critically relies on linearity condition

## Principal Hessian Direction (PHD)

- ▶ PHD estimate of  $\mathcal{S}_{E(Y|\mathbf{x})}$ : non-zero eigen-space of  $\Lambda_{\text{phd}} = E [\{Y - E(Y)\} \mathbf{x}\mathbf{x}^T]$
- ▶ Procedure:
  - ▶ Center  $Y$
  - ▶ Form sample average of  $\{Y - \hat{E}(Y)\} \mathbf{x}\mathbf{x}^T$ .
  - ▶ Extremely simple computation, finds all directions in  $\mathcal{S}_{E(Y|\mathbf{x})}$ .
- ▶ PHD requirement: critically relies on linearity condition and constant variance condition.

# Notes on Inverse Regression Estimators

- ▶ Elegant and mysterious
- ▶ Issues in tuning parameter selection: number of slices or bandwidth
- ▶ More and extended inverse regression estimators available both for  $\mathcal{S}_{Y|\mathbf{x}}$  and  $\mathcal{S}_{E(Y|\mathbf{x})}$
- ▶ Linearity/constant variance conditions bring some constraints
- ▶ Transformation often applied to reach these conditions
- ▶ Further assuming parametric model on  $\mathbf{x} \mid Y$  bypasses high dimensional issue, but assumption too strong

# Nonparametric Methods

- ▶ General idea of nonparametric methods: Minimize a criterion that describes the goodness-of-fit of the model.
- ▶ MAVE and dMAVE
  - ▶ Least square criterion
  - ▶ Nonparametric estimation involved: Local linear estimator
- ▶ Requirement: All components in  $\mathbf{x}$  are continuous

## Minimum Average Variance Estimation (MAVE)

- Local linear estimator for  $m(\mathbf{z})$  in  $Y = m(\mathbf{z}) + \epsilon$

$$\min_{a, \mathbf{b}} \sum_{i=1}^n \{ Y_i - a - \mathbf{b}^T (\mathbf{z}_i - \mathbf{z}) \}^2 K_h(\mathbf{z}_i - \mathbf{z})$$

- $\mathcal{S}_{E(Y|\mathbf{x})}$  model  $Y = m(\beta^T \mathbf{x}) + \epsilon$
- MAVE for  $\mathcal{S}_{E(Y|\mathbf{x})}$ 
  - let  $w_{ij} = K_h(\beta^T \mathbf{x}_i - \beta^T \mathbf{x}_j)$ .
  - Set  $\mathbf{z} = \beta^T \mathbf{x}_j$ 's.

$$\min \sum_{j=1}^n \sum_{i=1}^n \{ Y_i - a_j - \mathbf{b}_j^T (\beta^T \mathbf{x}_i - \beta^T \mathbf{x}_j) \}^2 w_{ij}$$

- Other criterion possible
- Other nonparametric estimator possible

## Density Based MAVE (dMAVE)

- dMAVE for  $\mathcal{S}_{Y|\mathbf{x}} : f_{Y|\mathbf{x}}(y, \mathbf{x}) = f_{Y|\beta^T \mathbf{x}}(y, \beta^T \mathbf{x})$ 
  - View it as  $I(Y = y) = m(\beta^T \mathbf{x}) + \epsilon$  for all  $y$  values.
  - MAVE yields

$$\min \sum_{j=1}^n \sum_{i=1}^n \{K_b(Y_i - y) - a_j - \mathbf{b}_j^T (\beta^T \mathbf{x}_i - \beta^T \mathbf{x}_j)\}^2 w_{ij}$$

- Select  $y = Y_1, \dots, Y_n$

$$\min \sum_{k=1}^n \sum_{j=1}^n \sum_{i=1}^n \{K_b(Y_i - Y_k) - a_{jk} - \mathbf{b}_{jk}^T (\beta^T \mathbf{x}_i - \beta^T \mathbf{x}_j)\}^2 w_{ij}$$

- Other criterion possible
- Other nonparametric estimator possible

# Notes on Nonparametric Estimators

- ▶ Intuitive and computationally intensive
- ▶ Tuning parameter selected via crossvalidation
- ▶ More and extended nonparametric estimators available both for  $S_{Y|x}$  and  $S_{E(Y|x)}$
- ▶ Do not need linearity/constant variance conditions
- ▶ Need continuity of each covariate. Could be relaxed with more computational complexity
- ▶ Often yields better performance than inverse regression based methods

# Semiparametric Methods

- $\mathcal{S}_{Y|\mathbf{x}}$  model is a semiparametric model:

$$f_{X,Y}(\mathbf{x}, Y; \beta, \eta) = \eta_1(\mathbf{x})\eta_2(Y, \beta^T \mathbf{x})$$

$\beta$ : of interest;  $\eta_1, \eta_2$ : nuisance

- Regular Asymptotically Linear (RAL) estimator:

$$n^{1/2}(\hat{\beta} - \beta) = n^{-1/2} \sum_{i=1}^n \phi(\mathbf{O}_i) + o_p(1)$$

$\phi(\mathbf{O}_i)$ : the  $i$ -th influence function.

- Estimator:  $\sum_{i=1}^n \phi(\mathbf{O}_i, \hat{\beta}) = 0, \text{var}(\sqrt{n}\hat{\beta}) = E\{\phi(\mathbf{O})\phi(\mathbf{O})^T\}$
- Find influence function  $\implies$  find RAL estimator

# The Geometry

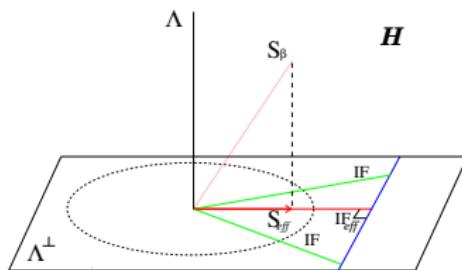
- ▶ Hilbert space

$$\begin{aligned}\mathcal{H}(\mathbf{O}) &= [\mathbf{f}(\mathbf{O}) : E\{\mathbf{f}(\mathbf{O})\} = 0, \text{var}\{\mathbf{f}(\mathbf{O})\} \text{ finite}] \\ <\mathbf{f}(\mathbf{O}), \mathbf{g}(\mathbf{O})> &= \text{cov}\{\mathbf{f}(\mathbf{O}), \mathbf{g}(\mathbf{O})\} = E\{\mathbf{f}(\mathbf{O})^T \mathbf{g}(\mathbf{O})\}\end{aligned}$$

- ▶  $\Lambda$  for parametric model:  $p(\mathbf{O}; \beta, \gamma)$ 
  - ▶ Nuisance score function:  $\mathbf{S}_\gamma(\mathbf{O}) = \partial \log p(\mathbf{O}; \beta, \gamma) / \partial \gamma$
  - ▶ Nuisance tangent space  $\Lambda_\gamma \subset \mathcal{H}$ : a linear space spanned by nuisance score function  $\mathbf{S}_\gamma$
- ▶  $\Lambda$  for semiparametric model:  $p(\mathbf{O}; \beta, \eta)$ 
  - ▶ Parametric submodel  $p(\mathbf{O}; \beta, \gamma)$
  - ▶ Nuisance tangent space  $\Lambda$ : span of  $\mathbf{S}_\gamma$  for all  $\gamma$ 's
- ▶  $\Lambda^\perp$ : Orthogonal complement of  $\Lambda$

# The Influence Functions

- ▶ Two requirements for influence function:
  1.  $\phi(\mathbf{O}) \in \Lambda^\perp$
  2.  $E\{\phi(\mathbf{O})\mathbf{S}_\beta^T(\mathbf{O})\} = \mathbf{I}$ ,  
 $\mathbf{S}_\beta(\mathbf{O}) = \partial \log p(\mathbf{O}) / \partial \beta$  is the score vector
- ▶ The spaces and influence functions visually



- ▶ Efficient influence function:

$$\phi_{\text{eff}}(\mathbf{O}) = E\{\mathbf{S}_{\text{eff}}(\mathbf{O})\mathbf{S}_{\text{eff}}(\mathbf{O})^T\}^{-1}\mathbf{S}_{\text{eff}}(\mathbf{O})$$

## $\Lambda^\perp$ for $\mathcal{S}_{Y|\mathbf{x}}$ and Many Possibilities

- ▶  $\Lambda = \Lambda_1 \oplus \Lambda_2$ , where

$$\Lambda_1 = \{\mathbf{f}(\mathbf{x}) : \forall \mathbf{f} \text{ such that } E(\mathbf{f}) = \mathbf{0}\}$$

$$\Lambda_2 = \{\mathbf{f}(Y, \beta^T \mathbf{x}) : \forall \mathbf{f} \text{ such that } E(\mathbf{f} | \mathbf{x}) = E(\mathbf{f} | \beta^T \mathbf{x}) = \mathbf{0}\}.$$

- ▶  $\Lambda^\perp = \{\mathbf{f}(Y, \mathbf{x}) - E(\mathbf{f} | \beta^T \mathbf{x}, Y) : E(\mathbf{f} | \mathbf{x}) = E(\mathbf{f} | \beta^T \mathbf{x}) \forall \mathbf{f}\}.$

- ▶  $\Lambda^\perp$  provides many candidate influence functions  
(unnormalized)

- ▶ For any  $\mathbf{g}, \mathbf{a}$ ,

$$\{\mathbf{g}(Y, \beta^T \mathbf{x}) - E(\mathbf{g} | \beta^T \mathbf{x})\} \{\alpha(\mathbf{x}) - E(\alpha | \beta^T \mathbf{x})\}$$

- ▶ For any  $\mathbf{g}_i, \mathbf{a}_i$ ,

$$\sum_{i=1}^k \{\mathbf{g}_i(Y, \beta^T \mathbf{x}) - E(\mathbf{g}_i | \beta^T \mathbf{x})\} \{\alpha_i(\mathbf{x}) - E(\alpha_i | \beta^T \mathbf{x})\}$$

# Double Centering Form and Double Robustness

- ▶ Original form

$$\{\mathbf{g}(Y, \beta^T \mathbf{x}) - E(\mathbf{g} | \beta^T \mathbf{x})\} \{\alpha(\mathbf{x}) - E(\alpha | \beta^T \mathbf{x})\}$$

- ▶ Mis-specify  $E(\alpha | \beta^T \mathbf{x})$  does not cause inconsistency.

$$E[\{\mathbf{g}(Y, \beta^T \mathbf{x}) - E(\mathbf{g} | \beta^T \mathbf{x})\} \{\alpha(\mathbf{x}) - \mathbf{h}(\beta^T \mathbf{x})\}] = \mathbf{0}$$

- ▶ Mis-specify  $E(\mathbf{g} | \beta^T \mathbf{x})$  does not cause inconsistency.

$$E[\{\mathbf{g}(Y, \beta^T \mathbf{x}) - \mathbf{h}(\beta^T \mathbf{x})\} \{\alpha(\mathbf{x}) - E(\alpha | \beta^T \mathbf{x})\}] = \mathbf{0}$$

This is extensively exploited in inverse regression type estimators.

- ▶ Cannot mis-specify both, but can estimate one or both nonparametrically.

## Linear Algebra Tools: Two Lemmas

- ▶ Lemma 1: Assume  $\Lambda$  is a  $p \times p$  symmetric matrix of rank  $d$ . If and only if  $\beta$  satisfies

$$\Lambda - \mathbf{P}\Lambda\mathbf{P} = \mathbf{0},$$

then the span of the columns in  $\beta$  is the eigen-space of  $\Lambda$  corresponding to the  $d$  nonzero eigenvalues.

- ▶ Lemma 2: Assume  $\Lambda$  is a  $p \times p$  symmetric non-negative definite matrix of rank  $d$ . If and only if  $\beta$  satisfies

$$\mathbf{Q}\Lambda\mathbf{Q} = \mathbf{0},$$

then the span of the columns in  $\beta$  is the eigen-space of  $\Lambda$  corresponding to the  $d$  nonzero eigenvalues.

# SIR and Semi-SIR

- ▶ SIR
  - ▶ SIR estimate of  $\mathcal{S}_{Y|\mathbf{x}}$ : non-zero eigen-space of  $\Lambda_{\text{sir}} = \text{cov}\{E(\mathbf{x} | Y)\}$
  - ▶ SIR requirement: linearity condition
- ▶ Rederive SIR from semiparametric method
  - ▶ Set  $\mathbf{g}(Y, \beta^T \mathbf{x}) = E(\mathbf{x} | Y)$ ,  $\alpha(\mathbf{x}) = \mathbf{x}^T$ .
  - ▶  $E\{\alpha(\mathbf{x}) | \beta^T \mathbf{x}\} = \mathbf{x}^T \mathbf{P}$ . Mis-specify  $E(\mathbf{g} | \beta^T \mathbf{x}) = 0$ .
  - ▶ Result in  $\Lambda_{\text{sir}} \mathbf{Q} = \mathbf{0}$ .
  - ▶ Lemma 2 ensures it is exactly SIR.
- ▶ Improve SIR by relaxing linearity condition, Semi-SIR:

$$E(\mathbf{x} | Y) \left\{ \mathbf{x} - \hat{E}(\mathbf{x} | \beta^T \mathbf{x}) \right\}^T$$

# SAVE and Semi-SAVE

## ► SAVE

- SAVE estimate of  $S_{Y|x}$ : non-zero eigen-space of

$$\Lambda_{\text{save}} = E \left[ \{\mathbf{I}_p - \text{cov}(\mathbf{x} | Y)\}^2 \right]$$

- SAVE restriction: linearity and constant variance condition

## ► Rederive SAVE from semiparametric method

- $\mathbf{g}_1(Y, \beta^T \mathbf{x}) = \mathbf{I}_p - \text{cov}(\mathbf{x} | Y)$ ,  $\mathbf{g}_2(Y, \beta^T \mathbf{x}) = \mathbf{g}_1 E(\mathbf{x} | Y)$ ,

$$\alpha_1(\mathbf{x}) = -\mathbf{x}\{\mathbf{x} - E(\mathbf{x} | \beta^T \mathbf{x})\}^T, \alpha_2(\mathbf{x}) = \mathbf{x}^T.$$

- Calculate  $E(\alpha_1 | \beta^T \mathbf{x})$ ,  $E(\alpha_2 | \beta^T \mathbf{x})$ .

Mis-specify  $E(\mathbf{g}_1 | \beta^T \mathbf{x}) = E(\mathbf{g}_2 | \beta^T \mathbf{x}) = \mathbf{0}$ .

- Result in  $\Lambda_{\text{save}} \mathbf{Q} = \mathbf{0}$ .

- Lemma 2 ensures it is exactly SAVE.

## ► Improve SAVE by relaxing both conditions, Semi-SAVE:

$$\{\mathbf{I}_p - \text{cov}(\mathbf{x} | Y)\} \left[ \{\mathbf{x} - E(\mathbf{x} | Y)\} \left\{ \mathbf{x} - \widehat{E}(\mathbf{x} | \beta^T \mathbf{x}) \right\}^T - \widehat{\text{cov}}(\mathbf{x} | \beta^T \mathbf{x}) \right]$$

# DR and Semi-DR

## ► DR

- ▶ Non-zero eigen-space of  $\Lambda_{\text{dr}} = E[\{2\mathbf{I}_p - \mathbf{A}(Y, \tilde{Y})\}^2]$ ,  
 $\mathbf{A}(Y, \tilde{Y}) = E\{(\mathbf{x} - \tilde{\mathbf{x}})(\mathbf{x} - \tilde{\mathbf{x}})^T \mid Y, \tilde{Y}\}$ ,  $(\tilde{\mathbf{x}}, \tilde{Y}) \perp\!\!\!\perp (\mathbf{x}, Y)$ .
- ▶ DR requirement: linearity and constant variance condition
- ▶ Rederive DR from semiparametric method
  - ▶  $\mathbf{g}_1(Y, \beta^T \mathbf{x}) = \mathbf{I}_p - E(\mathbf{x}\mathbf{x}^T \mid Y)$ ,
  - ▶  $\mathbf{g}_2(Y, \beta^T \mathbf{x}) = E\{E(\mathbf{x} \mid Y)E(\mathbf{x}^T \mid Y)\} E(\mathbf{x} \mid Y)$
  - ▶  $\mathbf{g}_3(Y, \beta^T \mathbf{x}) = E\{E(\mathbf{x}^T \mid Y)E(\mathbf{x} \mid Y)\} E(\mathbf{x} \mid Y)$
  - ▶  $\alpha_1(\mathbf{x}) = -\mathbf{x}\{\mathbf{x} - E(\mathbf{x} \mid \beta^T \mathbf{x})\}^T$ ,  $\alpha_2(\mathbf{x}) = \alpha_3(\mathbf{x}) = \mathbf{x}^T$ .
  - ▶ Calculate  $E(\alpha_i \mid \beta^T \mathbf{x})$ . Mis-specify  $E(\mathbf{g}_i \mid \beta^T \mathbf{x}) = \mathbf{0}$ .
  - ▶ Result in  $\Lambda_{\text{dr}} \mathbf{Q} = \mathbf{0}$ .
  - ▶ Lemma 2 ensures it is exactly DR.
- ▶ Improve DR by relaxing both conditions, Semi-DR

$$\begin{aligned} & \{\mathbf{I}_p - E(\mathbf{x}\mathbf{x}^T \mid Y)\} \left\{ -\mathbf{x}\mathbf{x}^T + \mathbf{x}\widehat{E}(\mathbf{x}^T \mid \beta^T \mathbf{x}) + \widehat{\text{cov}}(\mathbf{x} \mid \beta^T \mathbf{x}) \right\} \\ & + E\{E(\mathbf{x} \mid Y)E(\mathbf{x}^T \mid Y)\} E(\mathbf{x} \mid Y) \left\{ \mathbf{x}^T - \widehat{E}(\mathbf{x}^T \mid \beta^T \mathbf{x}) \right\} \\ & + E\{E(\mathbf{x}^T \mid Y)E(\mathbf{x} \mid Y)\} E(\mathbf{x} \mid Y) \left\{ \mathbf{x}^T - \widehat{E}(\mathbf{x}^T \mid \beta^T \mathbf{x}) \right\} \end{aligned}$$

## Link to Non-elliptical SIR

- ▶ nSIR
  - ▶ nSIR estimate of  $\mathcal{S}_{Y|\mathbf{x}}$ : minimizing
$$E \left( \| E(\mathbf{x} | Y) - E \{ E(\mathbf{x} | \boldsymbol{\beta}^T \mathbf{x}) | Y \} \|^2 \right)$$
  - ▶ nSIR restriction: polynomial condition and constant variance condition
- ▶ Rederive nSIR from semiparametric method:
  - ▶  $\mathbf{a}(\mathbf{x}) = \mathbf{x} - E(\mathbf{x} | \boldsymbol{\beta}^T \mathbf{x})$
  - ▶  $\mathbf{g}(Y, \boldsymbol{\beta}^T \mathbf{x}) = E \left[ \partial E(\mathbf{x} | \boldsymbol{\beta}^T \mathbf{x}) / \partial \{ \text{vec}(\boldsymbol{\beta}) \}^T | Y \right]$
  - ▶  $E(\mathbf{a} | \boldsymbol{\beta}^T \mathbf{x}) = \mathbf{0}$  guaranteed. Mis-specify  $E(\mathbf{g} | \boldsymbol{\beta}^T \mathbf{x}) = \mathbf{0}$ .
  - ▶ Result in nSIR.
- ▶ Improve nSIR by relaxing both conditions

# Link to Non-elliptical SSAVE

- ▶ nSAVE

- ▶ nSAVE estimate of  $S_{Y|x}$ : minimizing

$$E \left( \left\| \mathbf{I}_p - \text{cov}(\mathbf{x} | Y) - \text{cov} \{ E(\mathbf{x} | \boldsymbol{\beta}^T \mathbf{x}) \} + \text{cov} \{ E(\mathbf{x} | \boldsymbol{\beta}^T \mathbf{x}) | Y \} \right\|^2 \right)$$

- ▶ Requirement: polynomial condition and constant variance condition

- ▶ Rederive nSAVE from semiparametric method

- ▶  $\mathbf{g}(Y, \boldsymbol{\beta}^T \mathbf{x}) = \text{vec} \{ \mathbf{I}_p - \text{cov}(\mathbf{x} | Y) - \text{cov}(\mathbf{u}) + \text{cov}(\mathbf{u} | Y) \}^T$

$$\boldsymbol{\alpha}(\mathbf{x}) = \frac{\partial \text{vec}(\mathbf{u} \mathbf{u}^T)}{\partial \text{vec}(\boldsymbol{\beta})^T} - \mathbf{u} \otimes \frac{\partial \mathbf{u}}{\partial \text{vec}(\boldsymbol{\beta})^T} - \frac{\partial \mathbf{u}}{\partial \text{vec}(\boldsymbol{\beta})^T} \otimes \mathbf{u}$$

- ▶  $E(\mathbf{g} | \boldsymbol{\beta}^T \mathbf{x}) = \mathbf{0}$  guaranteed. Mis-specify  $E(\boldsymbol{\alpha} | \boldsymbol{\beta}^T \mathbf{x}) = \mathbf{0}$
  - ▶ Result in nSAVE.

- ▶ Improve nSAVE by relaxing both conditions

## Semiparametric Estimation of $\mathcal{S}_{E(Y|\mathbf{x})}$

- ▶ Model:  $E(Y | \mathbf{x}) = E(Y | \beta^T \mathbf{x})$
- ▶ Equivalent form

$$Y = m(\beta^T \mathbf{x}) + \epsilon, E(\epsilon | \mathbf{x}) = 0$$

- ▶ pdf:  $\eta_1(\mathbf{x})\eta_2\{Y - m(\beta^T \mathbf{x}), \mathbf{x}\}$ , where  $\int \epsilon \eta_2(\epsilon, \mathbf{x}) d\mu(\epsilon) = 0$
- ▶  $\Lambda = \Lambda_1 \oplus \Lambda_2$ , where

$$\Lambda_1 = \{\mathbf{f}(\mathbf{x}) : \forall \mathbf{f} \text{ such that } E(\mathbf{f}) = \mathbf{0}\}$$

$$\Lambda_2 = \{\mathbf{f}(\epsilon, \mathbf{x}) : \forall \mathbf{f} \text{ such that } E(\epsilon \mathbf{f} | \mathbf{x}) = E(\mathbf{f} | \mathbf{x}) = \mathbf{0}\}.$$

- ▶  $\Lambda^\perp = [\{Y - E(Y | \beta^T \mathbf{x})\} \{\alpha(\mathbf{x}) - E(\alpha | \beta^T \mathbf{x})\} : \forall \alpha]$
- ▶ Fewer choices.
- ▶ Can still mis-specify either  $E(Y | \beta^T \mathbf{x})$  or  $E(\alpha | \beta^T \mathbf{x})$ .

# OLS and Semi-OLS

- ▶ OLS
  - ▶ OLS estimate of  $S_{E(Y|\mathbf{x})}$ : one subspace is  $\text{cov}(\mathbf{x}, Y)$
  - ▶ OLS requirement: linearity condition
- ▶ Rederive OLS from semiparametric method
  - ▶ Set  $\alpha(\mathbf{x}) = \mathbf{x}$ .
  - ▶  $E\{\alpha(\mathbf{x}) | \beta^T \mathbf{x}\} = \mathbf{P}\mathbf{x}$ . Mis-specify  $E(Y | \beta^T \mathbf{x}) = 0$
  - ▶ Obtain
$$\mathbf{0} = E(\mathbf{x}Y) - \mathbf{P}E(\mathbf{x}Y) = E(\mathbf{x}Y) - \beta(\beta^T \beta)^{-1} \beta^T E(\mathbf{x}Y)$$
This is exactly OLS.
- ▶ Improve OLS by relaxing linearity condition, Semi-OLS

$$Y \left\{ \mathbf{x} - \widehat{E}(\mathbf{x} | \beta^T \mathbf{x}) \right\}$$

# PHD and Semi-PHD

- ▶ PHD
  - ▶ PHD estimate of  $S_{E(Y|\mathbf{x})}$ : non-zero eigen-space of  $\Lambda_{\text{phd}} = E [\{Y - E(Y)\} \mathbf{x}\mathbf{x}^T]$
  - ▶ PHD requirement: linearity and constant variance condition
- ▶ Rederive PHD from semiparametric method
  - ▶ Set  $\alpha(\mathbf{x}) = \mathbf{x}\mathbf{x}^T$ .
  - ▶  $\alpha(\mathbf{x}) - E \{\alpha(\mathbf{x}) \mid \beta^T \mathbf{x}\} = \mathbf{x}\mathbf{x}^T - \mathbf{Q} - \mathbf{P}\mathbf{x}\mathbf{x}^T\mathbf{P}$   
Mis-specify  $E(Y \mid \beta^T \mathbf{x}) = E(Y)$ .
  - ▶ Obtain
$$E [\{Y - E(Y)\} (\mathbf{x}\mathbf{x}^T - \mathbf{P}\mathbf{x}\mathbf{x}^T\mathbf{P})] = \Lambda_{\text{phd}} - \mathbf{P}\Lambda_{\text{phd}}\mathbf{P} = \mathbf{0}$$
  - ▶ Lemma 1 ensures it is exactly PHD
- ▶ Improve PHD by relaxing both conditions, Semi-PHD

$$\{Y - E(Y)\} \left\{ \mathbf{x}\mathbf{x}^T - \widehat{E}(\mathbf{x}\mathbf{x}^T \mid \beta^T \mathbf{x}) \right\}$$

# Simulation

- ▶  $p = 12, d = 2, n = 200, 500$  simulations
- ▶ Four models:

$$\text{model 1 : } Y = (\mathbf{x}^T \boldsymbol{\beta}_1) / \left\{ 0.5 + (\mathbf{x}^T \boldsymbol{\beta}_2 + 1.5)^2 \right\} + 0.5\epsilon;$$

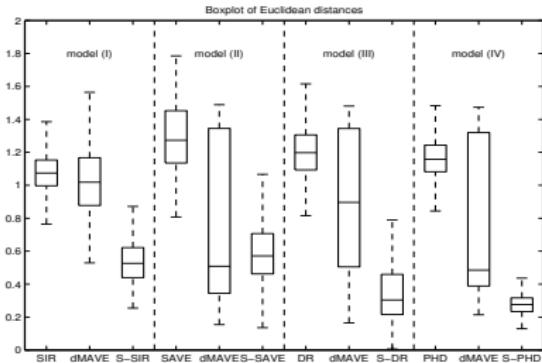
$$\text{model 2 : } Y = (\mathbf{x}^T \boldsymbol{\beta}_1)^2 + 2 |\mathbf{x}^T \boldsymbol{\beta}_2| + 0.1 |\mathbf{x}^T \boldsymbol{\beta}_2| \epsilon;$$

$$\text{model 3 : } Y = \exp(\mathbf{x}^T \boldsymbol{\beta}_1) + 2 (\mathbf{x}^T \boldsymbol{\beta}_2 + 1)^2 + |\mathbf{x}^T \boldsymbol{\beta}_1| \epsilon;$$

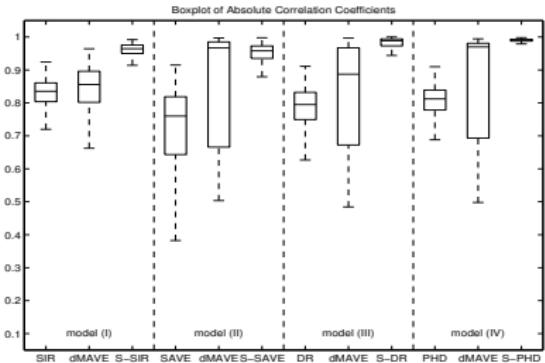
$$\text{model 4 : } Y = (\mathbf{x}^T \boldsymbol{\beta}_1)^2 + (\mathbf{x}^T \boldsymbol{\beta}_2)^2 + 0.5\epsilon,$$

- ▶ Two cases:
  - case 1:  $\mathbf{x}$  does not satisfy linearity or constant variance condition
  - case 2:  $\mathbf{x}$  satisfies both linearity and constant variance condition

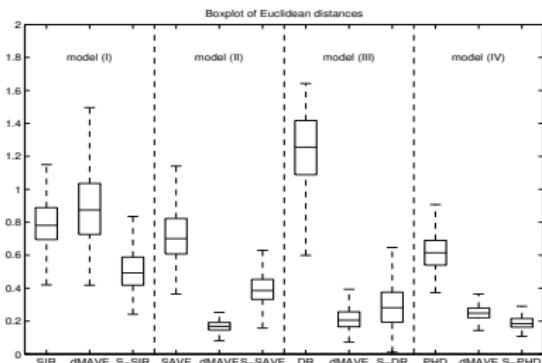
# Simulation Results and a Paradox



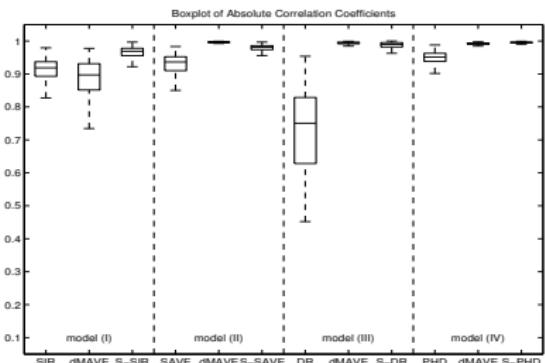
case 1: Euclidean distances



case 1: Canonical correlations



case 2: Euclidean distances



case 2: Canonical correlations

## Inference and Parameterization

- ▶ Inference on space estimation done via inference on parameter estimation
- ▶ Parameterization (up to a zero measure set)
  - ▶  $\beta$  is an orthogonal matrix: not unique
  - ▶ Only for  $d = 1$ :  $\beta_k = \sin(\theta_1) \dots \sin(\theta_{k-1}) \cos(\theta_k)$  for  $k < p$ ,  
 $\beta_p = \sin(\theta_1) \dots \sin(\theta_{p-2}) \sin(\theta_{p-1})$ : computation cumbersome
  - ▶ Let  $\beta_I \in \mathcal{R}^{(p-d) \times d}$  completely free,

$$\beta = \begin{pmatrix} \mathbf{I}_d \\ \beta_I \end{pmatrix}$$

- ▶ If  $\beta$  spans the dimension reduction space, so does  $\beta\mathbf{A}$ ,  
 $\forall \mathbf{A} \in \mathcal{R}^{d \times d}$  and invertable.
- ▶ A parameterization of a point on the Grassmann manifold  
(collection of  $d$ -dimensional subspaces of  $p$ -dimensional real space)
- ▶  $\text{vecl}(\mathbf{A})$ : vectorize lower  $(p - d) \times d$  submatrix of  $\mathbf{A}$

## Inference in Semiparametric Estimation

- ▶ Reducing the number of estimating equations to  $(p - d)d$  via GMM
- ▶ Standard asymptotic analysis to obtain large sample properties on  $\hat{\beta}_I$
- ▶ Typically, consistent estimation in weight matrix has no first order effect.
- ▶ Typically, if properly done, starting from a function in  $\Lambda^\perp$ , additional nonparametric estimation has no first order effect.

# Is the Paradox Real?

- ▶ Recall  $\{\mathbf{g}(Y, \beta^T \mathbf{x}) - E(\mathbf{g} | \beta^T \mathbf{x})\} \{\mathbf{a}(\mathbf{x}) - E(\mathbf{a} | \beta^T \mathbf{x})\}$
- ▶ Inverse regression methods solve

$$\sum_{i=1}^n \{\mathbf{g}(Y_i, \beta^T \mathbf{x}_i) - h(\beta^T \mathbf{x}_i)\} \{\mathbf{a}(\mathbf{x}_i) - E(\mathbf{a} | \beta^T \mathbf{x}_i)\} = \mathbf{0}$$

- ▶ Semiparametric estimators solve

$$\sum_{i=1}^n \{\mathbf{g}(Y_i, \beta^T \mathbf{x}_i) - h(\beta^T \mathbf{x}_i)\} \{\mathbf{a}(\mathbf{x}_i) - \hat{E}(\mathbf{a} | \beta^T \mathbf{x}_i)\} = \mathbf{0}$$

- ▶ It is real! Formal asymptotic analysis shows semiparametric estimator more efficient than inverse regression estimator.

# Intuitive Understanding

- ▶ General parametric models:  $\mathbf{a}(\mathbf{x}) = \mathbf{m}(\boldsymbol{\beta}^T \mathbf{x}, \boldsymbol{\alpha}) + \epsilon$   
Solve  $\sum_{i=1}^n \mathbf{A}(\boldsymbol{\beta}^T \mathbf{x}_i) \{ \mathbf{a}(\mathbf{x}_i) - \mathbf{m}(\boldsymbol{\beta}^T \mathbf{x}_i, \boldsymbol{\alpha}) \} = \mathbf{0}$  to obtain  $\hat{\boldsymbol{\alpha}}(\boldsymbol{\beta})$ .

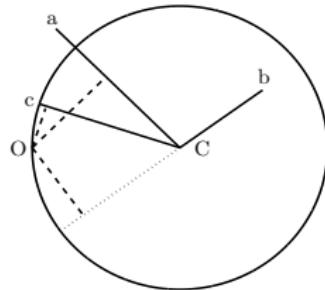
- ▶ Three estimators

$$(L) \quad \sum_{i=1}^n \mathbf{g}(Y_i, \boldsymbol{\beta}^T \mathbf{x}_i) \{ \mathbf{a}(\mathbf{x}_i) - \mathbf{m}(\boldsymbol{\beta}^T \mathbf{x}_i) \} = \mathbf{0},$$

$$(G) \quad \sum_{i=1}^n \mathbf{g}(Y_i, \boldsymbol{\beta}^T \mathbf{x}_i) \{ \mathbf{a}(\mathbf{x}_i) - \mathbf{m}(\boldsymbol{\beta}^T \mathbf{x}_i, \hat{\boldsymbol{\alpha}}(\boldsymbol{\beta})) \} = \mathbf{0},$$

$$(N) \quad \sum_{i=1}^n \mathbf{g}(Y_i, \boldsymbol{\beta}^T \mathbf{x}_i) \{ \mathbf{a}(\mathbf{x}_i) - \hat{\mathbf{m}}(\boldsymbol{\beta}^T \mathbf{x}_i) \} = \mathbf{0}.$$

- ▶  $\text{var}(L) = \text{var}(N) + \text{var}(C)$ ,  $\text{var}(G) = \text{var}(N) + \text{var}(a/b/c)$ . It is not true that  $L > G > N$ .



# Under OWLS

- ▶ OWLS:

$$\sum_{i=1}^n \mathbf{m}'(\beta^T \mathbf{x}_i, \alpha) \mathbf{V}^{-1}(\beta^T \mathbf{x}_i) \{\mathbf{a}(\mathbf{x}_i) - \mathbf{m}(\beta^T \mathbf{x}_i, \alpha)\} = \mathbf{0}$$

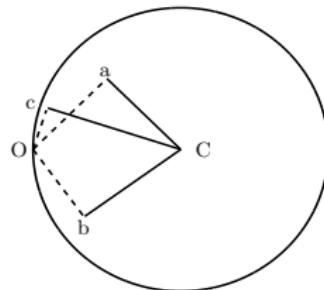
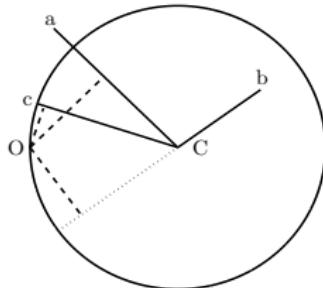
- ▶ Three estimators

$$(L) \quad \sum_{i=1}^n \mathbf{g}(Y_i, \beta^T \mathbf{x}_i) \{\mathbf{a}(\mathbf{x}_i) - \mathbf{m}(\beta^T \mathbf{x}_i)\} = \mathbf{0},$$

$$(OG) \quad \sum_{i=1}^n \mathbf{g}(Y_i, \beta^T \mathbf{x}_i) \{\mathbf{a}(\mathbf{x}_i) - \mathbf{m}(\beta^T \mathbf{x}_i, \hat{\alpha}(\beta))\} = \mathbf{0},$$

$$(N) \quad \sum_{i=1}^n \mathbf{g}(Y_i, \beta^T \mathbf{x}_i) \{\mathbf{a}(\mathbf{x}_i) - \hat{\mathbf{m}}(\beta^T \mathbf{x}_i)\} = \mathbf{0}.$$

- ▶  $\text{var}(L) = \text{var}(N) + \text{var}(C)$ ,  $\text{var}(OG) = \text{var}(N) + \text{var}(a/b/c)$ .  
Always  $L > OG > N$ .



## Nested Model and Best Estimator

- ▶ For nested models,  $a/b/c$  becomes distance to a line, a plain, a hyperplain...
- ▶ As dimension of the hyperplain increases, the distance decreases.
- ▶ Nonparametric kernel estimation is the extreme case of OWLS.
- ▶ This is the geometric intuition of why ignoring linearity/constant variance leads to better estimator.

# A Fair Evaluation of Linearity Condition

- ▶ Inverse regression methods rely on Linearity/Constant Variance condition to be consistent, the price is inflation of variability
- ▶ Linearity/Constant Variance condition does not bring loss if not overused
- ▶ Linearity/Constant Variance condition brings no gain if not taken into account at beginning
- ▶ Optimal use of these conditions improves efficiency compared to without these conditions. Efficient estimator forms very complex.

## Efficient Estimation of $\mathcal{S}_{Y|\mathbf{x}}$

- ▶ Score function:  $\mathbf{S}_\beta = \text{vecl}\{\mathbf{x} \partial \log \eta_2(Y, \beta^T \mathbf{x}) / \partial(\mathbf{x}^T \beta)\}$
- ▶ Efficient score

$$\mathbf{S}_{\text{eff}}(Y, \mathbf{x}, \beta^T \mathbf{x}, \eta_2, E) = \text{vecl} \left[ \{\mathbf{x} - E(\mathbf{x} \mid \beta^T \mathbf{x})\} \frac{\partial \log \eta_2(Y, \beta^T \mathbf{x})}{\partial(\mathbf{x}^T \beta)} \right]$$

- ▶ Estimate  $\eta_2, \partial \eta_2 / \partial(\mathbf{x}^T \beta), E(\cdot \mid \beta^T \mathbf{x})$ : all  $d$ -dimensional problems
- ▶ Efficient estimation

$$\sum_{i=1}^n \mathbf{S}_{\text{eff}}(Y_i, \mathbf{x}_i, \beta^T \mathbf{x}_i, \hat{\eta}_2, \hat{\eta}'_2, w, \hat{E}) = \mathbf{0}$$

- ▶  $\sqrt{n} \text{vecl}(\hat{\beta} - \beta) \rightarrow N[\mathbf{0}, \{\mathbf{S}_{\text{eff}}^{\otimes 2}(Y, \mathbf{x}, \beta^T \mathbf{x}, \eta_2, E)\}^{-1}]$ .

## Efficient Estimation of $\mathcal{S}_{E(Y|\mathbf{x})}$

- ▶ Score function

$$\mathbf{S}_\beta = -\text{vecl} \left\{ \frac{\mathbf{x} \eta'_{2\epsilon}(\epsilon, \mathbf{x}) \partial m(\beta^T \mathbf{x})}{\eta_2(\epsilon, \mathbf{x}) \partial (\mathbf{x}^T \beta)} \right\}$$

- ▶ Efficient estimation of  $\mathcal{S}_{E(Y|\mathbf{x})}$ . Let  $w(\mathbf{x}) = 1/E(\epsilon^2 | \mathbf{x})$ .

$$\begin{aligned}\mathbf{S}_{\text{eff}}(Y, \mathbf{x}, \beta^T \mathbf{x}, m, E, w) \\= \text{vecl} \left( \epsilon w(\mathbf{x}) \left[ \mathbf{x} - \frac{E\{\mathbf{x}w(\mathbf{x}) | \beta^T \mathbf{x}\}}{E\{w(\mathbf{x}) | \beta^T \mathbf{x}\}} \right] \frac{\partial m(\beta^T \mathbf{x})}{\partial (\mathbf{x}^T \beta)} \right)\end{aligned}$$

- ▶ Estimate  $m, m', E(\cdot | \beta^T \mathbf{x})$ : all  $d$ -dimensional problems
- ▶ Estimate  $w$ :  $p$ -dimensional problem, consistency suffices
- ▶ Efficient estimation

$$\sum_{i=1}^n \mathbf{S}_{\text{eff}}(Y_i, \mathbf{x}_i, \beta^T \mathbf{x}_i, \hat{m}, \hat{m}', \hat{E}(\cdot | \beta^T \mathbf{x}), \hat{w}) = \mathbf{0}$$

- ▶  $\sqrt{n} \text{vecl}(\hat{\beta} - \beta) \rightarrow N[\mathbf{0}, \{\mathbf{S}_{\text{eff}}^{\otimes 2}(Y, \mathbf{x}, \beta^T \mathbf{x}, m, E, w)\}^{-1}]$

# Simulation $\mathcal{S}_{Y|x}$

		$\beta_1$ 1.3	$\beta_2$ -1.3	$\beta_3$ 1	$\beta_4$ -5	$\beta_5$ .5	$\beta_6$ -.5
Ora	ave	1.2978	-1.3036	1.0049	-0.4985	0.5033	-0.4943
	std	0.1221	0.1477	0.1505	0.1169	0.0966	0.1049
	$\widehat{\text{std}}$	0.1264	0.1510	0.1527	0.1212	0.0983	0.1052
	95%	0.9510	0.9540	0.9440	0.9540	0.9520	0.9450
Eff	ave	1.2980	-1.3046	1.0064	-0.4990	0.5040	-0.4936
	std	0.1280	0.1546	0.1567	0.1221	0.1000	0.1075
	$\widehat{\text{std}}$	0.1317	0.1588	0.1602	0.1264	0.1011	0.1084
	95%	0.9480	0.9380	0.9380	0.9440	0.9480	0.9510
Loc	ave	1.3052	-1.2629	0.9687	-0.4988	0.5023	-0.4897
	std	0.1478	0.1736	0.1715	0.1393	0.1069	0.1153
dMA	ave	1.2599	-1.2933	1.0014	-0.4763	0.4984	-0.4935
	std	0.1932	0.1427	0.1550	0.1701	0.1368	0.1378
SIR	ave	1.3881	-1.1930	0.9261	-0.5968	0.4793	-0.4724
	std	0.1696	0.1522	0.1414	0.1489	0.0976	0.0995
DR	ave	0.9935	-0.2217	0.1930	-0.6863	0.1245	-0.1071
	std	0.6567	1.2305	1.0107	0.6411	0.3069	0.2999

# Simulation $\mathcal{S}_{E(Y|\mathbf{x})}$

		$\beta_1$ -1	$\beta_2$ -2	$\beta_3$ -0.5	$\beta_4$ -1	$\beta_5$ -2	$\beta_6$ -2	$\beta_7$ -1.5	$\beta_8$ -1
o	$\hat{\beta}$	-1.03	-2.08	-0.52	-1.04	-2.05	-2.05	-1.55	-1.02
r	$\sigma$	0.11	0.14	0.08	0.08	0.14	0.16	0.11	0.10
c	CI	95.6	93.0	93.9	94.2	94.9	94.2	93.6	94.4
I	$\hat{\beta}$	-1.04	-2.09	-0.52	-1.05	-2.06	-2.05	-1.55	-1.03
o	$\sigma$	0.11	0.14	0.08	0.08	0.15	0.17	0.11	0.10
1	CI	94.6	92.4	94.1	92.9	93.2	93.2	92.8	92.7
I	$\hat{\beta}$	-1.02	-2.05	-0.51	-1.03	-2.04	-2.04	-1.54	-1.02
o	$\sigma$	0.11	0.15	0.08	0.09	0.16	0.18	0.12	0.11
2	CI	94.9	94.0	94.8	94.1	94.1	94.6	93.6	92.8
I	$\hat{\beta}$	-1.03	-2.07	-0.51	-1.04	-2.05	-2.04	-1.54	-1.02
o	$\sigma$	0.11	0.14	0.08	0.08	0.14	0.16	0.11	0.10
3	CI	95.1	94.0	94.7	93.5	95.2	94.2	94.5	94.3
M	$\hat{\beta}$	-1.01	-1.88	-0.52	-1.06	-1.91	-1.86	-1.48	-1.01
A	$\sigma$	0.17	0.20	0.13	0.12	0.22	0.26	0.18	0.16
S	$\hat{\beta}$	-1.00	-1.92	-0.51	-1.03	-2.03	-1.94	-1.54	-1.04
P	$\sigma$	0.17	0.19	0.12	0.15	0.36	0.34	0.27	0.30

## Determine $d$ in Inverse Regression Based Methods

- ▶ Inverse regression based method procedure
  - ▶ Estimate  $\Lambda$  from data, get  $\widehat{\Lambda}$ .
  - ▶ Estimate the nonzero eigenspace of  $\widehat{\Lambda}$ .
- ▶  $d$  is the number of non-zero eigenvalues of the matrix  $\Lambda$ .
- ▶ Deciding  $d$  is to determine the number of non-zero eigenvalues of  $\Lambda$  based on  $\widehat{\Lambda}$ .

# Sequential Test

- ▶ Sequential test procedure:
  1.  $H_{01}$  :  $\Lambda$  has  $p - 1$  non-zero eigenvalues, if  $H_{01}$  passes,  
 $\hat{d} = p - 1$ , otherwise
  2.  $H_{02}$  :  $\Lambda$  has  $p - 2$  non-zero eigenvalues, if  $H_{02}$  passes,  
 $\hat{d} = p - 2$ , otherwise
  3. ...
- ▶ Sequential test problems:
  - ▶  $\hat{d}$  not consistent due to type-I errors allowed.
  - ▶ Cumulative type-I errors may be not ignorable.
  - ▶ When  $p$  large, computationally inefficient.
  - ▶ Heavily relies on normality of  $\hat{\Lambda}$ .
  - ▶ Needs  $\text{cov}(\hat{\Lambda})$ .
  - ▶ Done case-by-case for each inverse regression based estimator.

# BIC

- ▶ BIC procedure: The goodness-of-fit of a rank  $k$  matrix  $\Lambda_k$  to  $\widehat{\Lambda}$ ,  $f(\widehat{\Lambda}, \Lambda_k)$  is penalized by the rank  $k$ ,  $c(n, k)$ .
- ▶ General form  $\min_k f(\widehat{\Lambda}, \Lambda_k) + c(n, k)$
- ▶ BIC properties and problems:
  - ▶ It works with  $\widehat{\Lambda} + \mathbf{I}_p$  and  $\Lambda_k + \mathbf{I}_p$ , instead of  $\widehat{\Lambda}$  and  $\Lambda_k$  for technical reason.
  - ▶ Construction of goodness-of-fit  $f$  based on normality of  $\widehat{\Lambda}$ .
  - ▶ Designed only for SIR.
  - ▶ It is consistent as long as  $\widehat{\Lambda}$  is consistent.
  - ▶ It is sensitive to the amount of the penalty  $c(n, k)$ .

## Sparse Eigen-Decomposition (SED)

- Eigen decomposition problem is least square problem:

$$\min_{\lambda_i, \xi_i, \eta_i} \|\Lambda - \sum_{i=1}^p \lambda_i \xi_i \eta_i^T\|,$$

where  $\|\xi_i\| = \|\eta_i\| = 1, \lambda_1 \geq \dots \geq \lambda_p.$

- Equivalently minimize

$$\|\Lambda - \sum_{i=1}^p \lambda_i \xi_i \eta_i^T\|^2 = \|\Lambda\|^2 - 2 \sum_{i=1}^p \lambda_i \xi_i^T \Lambda \eta_i + \sum_{i=1}^p \lambda_i^2$$

to obtain  $\hat{\xi}_i(\lambda's), \hat{\eta}_i(\lambda's).$

- SED procedure (least square+lasso)

$$\min_{\lambda's} \|\Lambda - \sum_{i=1}^p \lambda_i \hat{\xi}_i(\lambda's) \hat{\eta}_i(\lambda's)^T\| + c(n) \sum_{i=1}^p \hat{w}_i |\lambda_i|$$

- SED property:  $\hat{d}$  and  $\hat{\lambda}'s$  consistent
- LARS applicable, computationally simple

# Bootstrap

- ▶ Intuition
  - ▶ Under true  $d$ , estimated space  $\mathbf{S}(d)$  approximates true space, has variability due to data randomness
  - ▶ Under  $k > d$ , estimated space  $\mathbf{S}(k)$  contains zero-eigendirection. Zero-eigendirection is not unique in population, the choice of zero-eigendirection in finite sample inflates variability
  - ▶ Under  $k < d$ , estimated space  $\mathbf{S}(k)$  is a subspace of  $\mathbf{S}(d)$ . The choice of the subspace in finite sample inflates variability.
- ▶ Variability of space estimation assessed via bootstrap
- ▶ Theoretical justification not available

# Determine $d$ in Nonparametric Methods

- ▶ Crossvalidation procedure

- ▶  $\mathcal{S}_{E(Y|\mathbf{x})}$

$$CV(k) = n^{-1} \sum_{i=1}^n \{Y_i - \hat{m}_{-i}(\hat{\beta}_{-i}^T \mathbf{x}_i)\}^2$$

- ▶  $\mathcal{S}_{Y|\mathbf{x}}$

$$CV(k) = n^{-1} \sum_{i=1}^n \sum_{j=1}^n \{K_b(Y_i - Y_j) - \hat{\eta}_{2,-i}(Y_j, \hat{\beta}_{-i}^T \mathbf{x}_i)\}^2$$

- ▶ Requires a goodness-of-fit criterion, leads to estimating quantities that could be avoided in other methods
- ▶ Computationally demanding

## Determine $d$ in Semiparametric Methods

- ▶ Estimating equation based estimation.
- ▶ Cannot be transformed to a problem of testing zero eigenvalues without linearity condition.
- ▶ No likelihood or other criterion computable, so AIC, BIC, CV do not apply.
- ▶ Quadratic form cannot replace likelihood as a criterion.
- ▶ Not a problem of zero coefficient: penalty does not apply.

# Bootstrap

- ▶ Intuition under true  $d$ 
  - ▶ Data based  $\hat{\mathbf{S}}(d)$  approximates true space  $\mathbf{S}(d)$
  - ▶ Bootstrap data based  $\hat{\mathbf{S}}_B(d)$  approximates true space  $\mathbf{S}(d)$
  - ▶  $\hat{\mathbf{S}}(d)$  and  $\hat{\mathbf{S}}_B(d)$  highly correlated
- ▶ Intuition under  $k > d$ 
  - ▶ Data based  $\hat{\mathbf{S}}(k)$  contains arbitrary direction out of  $\mathbf{S}(d)$
  - ▶ Bootstrap data based  $\hat{\mathbf{S}}_B(k)$  contains arbitrary direction out of  $\mathbf{S}(d)$
  - ▶ Arbitrary direction determined by data,  $\hat{\mathbf{S}}(k)$  and  $\hat{\mathbf{S}}_B(k)$  not as highly correlated
- ▶ Intuition under  $k < d$ 
  - ▶ Data based  $\hat{\mathbf{S}}(k)$  is arbitrary subspace of  $\mathbf{S}(d)$
  - ▶ Bootstrap data based  $\hat{\mathbf{S}}_B(k)$  is arbitrary subspace of  $\mathbf{S}(d)$
  - ▶ Arbitrary subspace determined by data,  $\hat{\mathbf{S}}(k)$  and  $\hat{\mathbf{S}}_B(k)$  not as highly correlated
- ▶ No theoretical proof for its validity

# VIC

- ▶ Key observation in semiparametric estimating equations
  1. If  $k < d$ , adding bad direction destroys estimating equation consistency under  $k + 1$  model.
  2. If  $k \geq d$ , adding bad direction does not destroy estimating equation consistency under  $k + 1$  model.
- ▶ VIC stands for Validated Information Criterion.
- ▶ Main idea: Extend a current  $\beta$  estimator by adding a column, examine if consistency is destroyed or not. Stop as soon as consistency does not get destroyed.
- ▶ We might get too lucky. Extend  $\hat{\beta}_{(k)}$  to  $\hat{\gamma}_{(k)}^1, \hat{\gamma}_{(k)}^2$ .

$$\begin{aligned} & \text{VIC}(k) \\ = & \frac{1}{2n} \left\{ \left\| \sum_{i=1}^n \mathbf{f}(\mathbf{O}_i, \hat{\gamma}_{(k)}^1) \right\|^2 + \left\| \sum_{i=1}^n \mathbf{f}(\mathbf{O}_i, \hat{\gamma}_{(k)}^2) \right\|^2 \right\} + pk\log(n) \end{aligned}$$

- ▶  $\hat{d} = \operatorname{argmin}_k \text{VIC}(k)$

## Properties of VIC

- ▶ Only need to solve for  $\beta$  in smaller ( $k$ ) model. No equation solving in larger ( $k + 1$ ) model.
- ▶ Theorem on selection consistency:  
Let  $\hat{d} = \operatorname{argmin}_k \text{VIC}(k)$ . Under regularity conditions, when  $n \rightarrow \infty$ ,  $\text{pr}(\hat{d} = d) \rightarrow 1$ .
- ▶ Simulation: repeated 1000 times. Conducted for  $n = 200, n = 400, p = 6, 10$  and  $d = 1, 2$ . Performed estimation after selecting structural dimension.

## Simulation Results: $d = 1$

Selection frequency by  $\text{VIC}(k)$  in %,  $p = 6$ .

	$\hat{d} = 1$	$\hat{d} = 2$	$\hat{d} = 3$	$\hat{d} = 4$
$n = 200$				
Semi-SIR	99.4	0.6	0.0	0.0
Semi-SAVE	94.1	5.2	0.6	0.1
Semi-DR	85.6	14.0	0.4	0.0
Semi-PHD	99.9	0.1	0.0	0.0
$n = 400$				
Semi-SIR	99.7	0.3	0.0	0.0
Semi-SAVE	94.4	5.5	0.1	0.0
Semi-DR	86.8	12.8	0.4	0.0
Semi-PHD	99.9	0.1	0.0	0.0

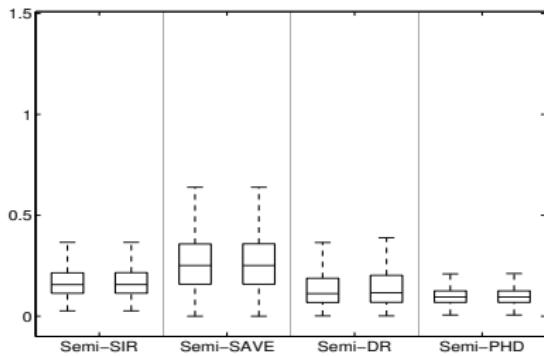
## Simulation Results: $d = 2$

Selection frequency by  $\text{VIC}(k)$  in %,  $p = 6$ .

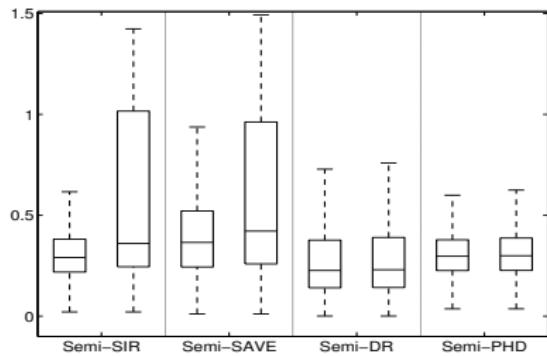
	$\hat{d} = 1$	$\hat{d} = 2$	$\hat{d} = 3$	$\hat{d} = 4$
$n = 200$				
Semi-SIR	31.0	68.3	0.6	0.1
Semi-SAVE	9.9	76.0	12.1	2.0
Semi-DR	1.4	96.6	1.7	0.3
Semi-PHD	2.6	97.2	0.2	0.0
$n = 400$				
Semi-SIR	4.6	94.9	0.5	0.0
Semi-SAVE	1.7	91.7	5.7	0.9
Semi-DR	1.6	98.1	0.2	0.1
Semi-PHD	0.0	99.8	0.2	0.0

# Estimation Results ( $p = 6$ , true and selected $d$ )

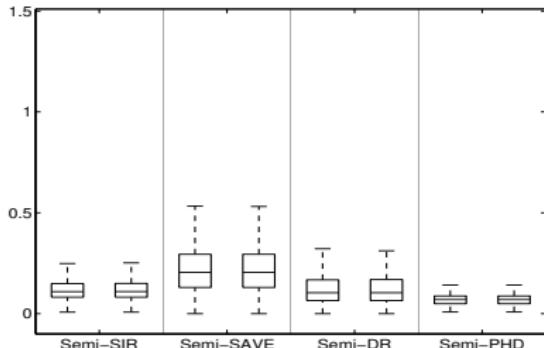
$d = 1, n = 200$



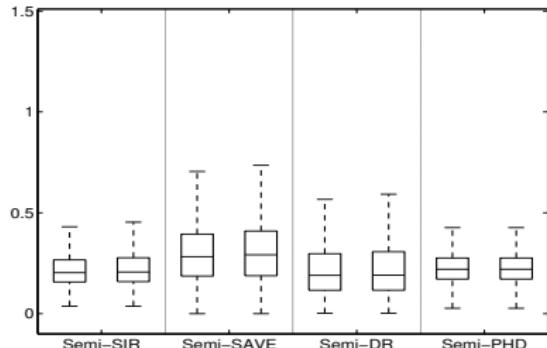
$d = 2, n = 200$



$d = 1, n = 400$



$d = 2, n = 400$



## Generality and Extension of VIC

- ▶ Supplement to AIC, BIC.
- ▶ Can be used in any situation when estimating equation exists, but criterion not available or hard to compute: No AIC, BIC, CV.
  - ▶ moment based methods
  - ▶ GEE
  - ▶ Martingale based methods
  - ▶ semiparametric methods
- ▶ Especially useful when model is nested, but smaller model is not result of having zero values in the larger model. No test, penalization.
- ▶ Especially useful when increasing model complexity implies great elevation of computation.

# When $p$ Is Very Large

- ▶ SIR valid for  $p = o(n^{1/4})$
- ▶ Cumulative slicing valid for  $p = o(n^{1/2})$
- ▶ Hypotheses
  - ▶ All current methods hold for  $p = o(n^{1/2})$
  - ▶ All current methods fail for  $p \neq o(n^{1/2})$
- ▶ Challenges and suggestions for  $p \gg n^{1/2}$ 
  - ▶  $\widehat{\Sigma}_x$  singular, normalization no longer apply
  - ▶  $\widehat{\Sigma}_x^{-1}$  not available, avoid it by partial least square idea
- ▶ Ultra large  $p$  requires additional assumptions
  - ▶ Screening
  - ▶ Assuming sparsity, performing penalizing
  - ▶ Introducing latent structure

# Extensions

- ▶ Extension to nonlinear case:  $\beta^T \mathbf{x} \implies \mathbf{u}(\mathbf{x}, \beta)$
- ▶ Extension to nonparametric case:  $\beta^T \mathbf{x} \implies \mathbf{u}(\mathbf{x})$
- ▶ Extension to functional dimension reduction
- ▶ Extension to two-way and multi-way dimension reduction (dimension folding)
- ▶ Extension to multivariate response
- ▶ Extension to incorporating additional constraints (envelope models)